

## NBER WORKING PAPER SERIES

### THE ROLE OF THEORY IN FIELD EXPERIMENTS

David Card  
Stefano DellaVigna  
Ulrike Malmendier

Working Paper 17047  
<http://www.nber.org/papers/w17047>

### NATIONAL BUREAU OF ECONOMIC RESEARCH

1050 Massachusetts Avenue  
Cambridge, MA 02138  
May 2011

We thank for helpful comments the editors, as well as Oriana Bandiera, Iwan Barankay, Glenn Harrison, Matthew Rabin, David Reiley, and participants in Berkeley, at the Wharton Conference on Field Experiment, and at the 2011 ASSA conference in Denver. We thank Ivan Balbuzanov, Xiaoyu Xia, and a very dedicated group of undergraduate students for excellent research assistance. The views expressed herein are those of the authors and do not necessarily reflect the views of the National Bureau of Economic Research.

NBER working papers are circulated for discussion and comment purposes. They have not been peer-reviewed or been subject to the review by the NBER Board of Directors that accompanies official NBER publications.

© 2011 by David Card, Stefano DellaVigna, and Ulrike Malmendier. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

The Role of Theory in Field Experiments  
David Card, Stefano DellaVigna, and Ulrike Malmendier  
NBER Working Paper No. 17047  
May 2011  
JEL No. C9,C93

### **ABSTRACT**

We propose a new classification of experiments that captures the extent to which the experimental design and analysis are linked to economic theory. We then use this system to classify all published field experiments in the five top economics journals from 1975 to 2010. We find that the vast majority of field experiments (68%) are Descriptive studies that lack any explicit model; 18% are Single Model studies that test a single model-based hypothesis; 6% are Competing Models studies that test competing model-based hypotheses; and 8% are Parameter Estimation studies that estimate structural parameters in a completely specified model. Using the same system to classify laboratory experiments published over the same period, we find that economic theory has played a more central role in the laboratory than in the field. Finally, we discuss in detail three sets of field experiments, on gift exchange, on charitable giving, and on negative income tax, that illustrate both the benefits and the potential costs of a tighter link between experimental design and theoretical underpinnings.

David Card  
Department of Economics  
549 Evans Hall, #3880  
University of California, Berkeley  
Berkeley, CA 94720-3880  
and NBER  
card@econ.berkeley.edu

Ulrike Malmendier  
Department of Economics  
549 Evans Hall # 3880  
University of California, Berkeley  
Berkeley, CA 94720-3880  
and NBER  
ulrike@econ.berkeley.edu

Stefano DellaVigna  
University of California, Berkeley  
Department of Economics  
549 Evans Hall #3880  
Berkeley, CA 94720-3880  
and NBER  
sdellavi@econ.berkeley.edu

When it comes to role of theory in their research, empirical microeconomists are torn. On the one hand, we devote a large fraction of our graduate instruction to models of consumer behavior and firm decision-making, and to the interactions that determine market equilibrium. On the other hand, it is not always obvious how these theories are relevant to empirical research. Outside the academy, policy-makers and business leaders often demand “basic facts” and simplified policy guidance with little or no concern for theoretical nuances.

How then do empirical economists negotiate between theory and “facts”? In this paper, we focus on the role of theory in the rapidly-growing area of field experiments. We take an empirical approach and quantify the role of theoretical modeling in all published field experiments in five top economics journals from 1975 to 2010. We propose a new classification of experimental studies that captures the extent to which the experimental design and analysis is linked to economic theory. Specifically, we distinguish between four classes of studies: *Descriptive* studies that lack any explicit model; *Single Model* studies that test a single model-based hypothesis; *Competing Models* studies that test competing model-based hypotheses; and *Parameter Estimation* studies that estimate structural parameters in a completely specified model. Applying the same classification to laboratory experiments published over the same period we conclude that theory has played a more central role in the laboratory than in field experiments. Finally, we discuss in detail three sets of field experiments that illustrate both the potential promise and pitfalls of a tighter link between experimental design and theoretical underpinnings.

### **Quantifying the Role of Theory in Field Experiments**

The use of “experimental” (i.e., random-assignment) designs came relatively late to economics.<sup>1</sup> Over the last 15 years, however, randomized experiments in field settings have proliferated, and in 2010 field experiments represented about 3 percent of the arti-

---

<sup>1</sup> According to Forsetlund, Chalmers, and Bjorndal (2007), the earliest documented use of randomization in the social sciences was a 1928 study of an intervention designed to reduce the rate at which college students were failing at Purdue University. In economics, we are unaware of any study using random assignment prior to the negative income tax experiments in the 1960s (see Greenberg and Shroder, 2004).

cles published in the top economics journals. The role of theory in such field experiments, as in other areas of applied economics, ranges from “almost none” to fully model-based investigations. However, there is a widespread perception that experimental studies, and particularly field-based random-assignment studies, are disproportionately “black box” evaluations that provide only limited evidence on theoretically-relevant mechanisms (for example, Deaton, 2010).

To assess the actual importance of theoretical modeling in field experiments, and compare the relative role of theory in field versus laboratory experiments, we collected data on the universe of experimental studies published in five leading economics journals – the *American Economic Review*, *Econometrica*, the *Journal of Political Economy*, the *Quarterly Journal of Economics*, and the *Review of Economic Studies* – over the 36-year period from 1975 to 2010. After excluding comments, notes, and articles in the annual *Papers and Proceedings Issue* of the *American Economic Review*, we identified all laboratory and field experiments among the remaining articles and classified the role of theory in these two sets of studies.

### *Defining Field Experiments*

A first issue that arises for our analysis is the delineation of what qualifies as an “experiment.” We restrict attention to studies based on the random assignment of a purposeful “treatment” or manipulation. We include studies where treatment is deterministically assigned in a way that can be viewed as equivalent to random, such as assigning every second name in a list, or choosing a permutation of potential subjects that optimizes the balance between treatment and control groups. Our definition includes government-funded social experiments, such as Moving To Opportunity, which provided vouchers for public housing recipients to move out of low-income neighborhoods, (Kling, Liebman, and Katz, 2007); smaller scale research projects like List’s (2003) study of sport card dealers; and randomizations induced by a firm for its own research or marketing purposes, like Nagin et al.’s (2002) study of the effects of monitoring on telephone solicitors.

However, our definition excludes many influential studies that are often viewed as “experimental” but lack randomly-assigned treatment and control groups. Bandiera,

Barankay, and Rasul's (2007, 2009) studies of bonus payments to farm managers, for example, use a "pre-post" design in which managers are first observed in one regime, and then in another. A similar non-random design is used by Chetty, Looney and Kroft (2009) to study the effect of including sales taxes in the posted prices displayed in grocery stores. We also exclude other studies that exploit random variation created for purposes other than the evaluation of treatment, like Angrist's (1990) study of the Vietnam draft lottery or Sacerdote's (2001) study of randomly-assigned college roommates.

By restricting attention to studies with random assignment of a purposeful manipulation, we do not mean to criticize papers that use non-randomized designs, or that rely on opportunistic randomization. Rather, we use these criteria to narrow our focus to studies that are closest in spirit to the *randomized clinical trials* used in medicine and other sciences. Advocates of randomized experimental studies often point to these trials as the gold standard for scientific evidence, despite the limitations emphasized by Heckman and Smith (1995) and Deaton (2010), for example.

We include papers that re-analyze data from previous experiments, provided that the study uses the original micro data, as in Lalonde's (1986) analysis of econometric methods for program evaluation. In the terminology of Harrison and List (2004) we include both "natural field experiments" in which the participants have no knowledge of being involved in an experiment and "framed field experiments" in which the participants are aware that they participate in an experiment.

### *Classification of the Role of Theory*

Within this universe of studies, we classify the role of economic theory using a four-way scheme that we believe captures the centrality of economic theory in a particular study. The four categories are: *Descriptive (D)* studies that lack any formally specified model; *Single Model (S)* studies that lay out a formal model and test one (or more) qualitative implications of the model; *Competing Models (C)* studies that lay out two or more alternative models with at least one contrasting qualitative implication and test between them on the basis of this implication; and *Parameter Estimation (P)* studies that specify a complete data generating process for (at least some subset of) the observed data and obtain estimates of structural parameters of the model.

To illustrate our classification system, Table 1 shows four examples from the recent literature that are broadly representative of the four classes. Miguel and Kremer’s (2004) study of a deworming treatment program in Kenya provides an interesting example of a *Descriptive* field experiment. The experimental treatment in this study contains several elements, including drug treatment and education, and was designed to affect a variety of outcomes, including infection rates, school attendance and educational achievement. The paper provides no formal model for the experimental program impacts, though it does discuss the expected effects on health and education outcomes as well as possible channels for these effects, including social spillovers.

*Single Model* experiments lay out a formal model of the experimental impact and then evaluate the predictions of this model against the null hypothesis of no difference between the treatment and control groups. To meet the definition of a “formal model” for this class we require at least one line of offset mathematical text. (We make no attempt to assess the logical completeness of the model specification). We exclude purely statistical models or algebraic summaries of the payoffs in laboratory experiments. An illustrative example is the paper by Nagin et al. (2002), which includes a simple but formally specified model that isolates the response of a key endogenous variable (the number of “questionable” calls claimed by a telephone sales associate) to a manipulation of interest (the monitoring rate of questionable calls). The qualitative prediction of the model is then tested by contrasting various treatment groups.

Although our requirement of a single equation of mathematical text provides an easily verified distinction between *Descriptive* and *Single-Model* studies, we readily concede that in some cases the line is arbitrary. Consider, for example, a field experiment designed to test an implication of a well-known model. In some cases a referees or editor will have asked the authors to remove the formal statement of the model from the paper, leading us to classify the paper as *Descriptive*. In other cases the formal statement remains, leading us to classify the paper as a *Single Model* study and inducing different classifications for papers which are equally informed by theory. Despite this issue, we believe that the presence of a mathematical statement of the model is a useful (if crude) indicator of the importance of economic theory in the paper. A formal statement of the model helps to clarify the underlying assumptions that the author is maintaining in the

study and the specific form of the model that the author is attempting to test in the empirical setting.<sup>2</sup>

A criticism of studies that focus on testing a single model is that they provide little guidance in the event that the model is rejected: Which of the assumptions does the data reject? Would alternative models have fared differently? A parallel criticism arises when the model is *not* rejected: Competing models may make the same prediction, so simple “one sided” tests do not distinguish between theories (Rabin, 2010). A text-book example of the latter problem is provided by Becker (1962), who notes that the finding of a downward-sloping demand curve cannot be construed as evidence of utility maximization, since demand curves will be downward-sloping even when agents choose randomly, as long the budget constraint is sometimes binding.

These concerns are partially addressed by *Competing Model* studies that lay out two or more competing models, with differing predictions for the response to a manipulation. The study by Fehr and Goette (2007), for example, compares a standard intertemporal labor supply model against an alternative model with reference-dependent preferences. The two models have similar predictions for the response of earnings to a short-term increase in the effective wage rate, but differing predictions for effort per hour: effort increases under the standard model, but decreases under reference dependence. The latter predictions provide the basis for a test between the models.

---

<sup>2</sup> A useful case to consider is the influential set of findings on the “disposition effect” – that is, on the propensity to sell stocks that are “winners” rather than “losers” compared to the purchase price. Odean (1999) uses a graph and an intuitive explanation to suggest that the phenomenon is explained by prospect theory. However, Barberis and Xiong (2009) show that once one actually writes down an explicit model of prospect theory, the disposition effect is not generally predicted by the model. In this case, the intuitive explanation had focused on the concavity and convexity of the value function, but had neglected the effect of the kink at the reference point.

**Table 1: Classification Examples**

| Study  | Description   | Classification              |
|--|---|-----------------------------|
| 1. E. Miguel and M. Kremer<br>"Worms: Identifying Impacts on Education and Health in the Presence of Treatment"<br><i>Econometrica</i> , 2004  | Evaluation of deworming treatment program in Kenya. School-level assignment. Treatment delayed at control groups.   | <b>Descriptive</b>          |
| 2. D. Nagin et al.<br>"Monitoring, Motivation, and Management: The Determinants of Opportunistic Behavior in a Field Experiment"<br><i>American Econ. Rev.</i> , 2002  | Random assignment of monitoring rate of call-center employees. Center-level assignment. Model of optimal cheating predicts greater cheating when monitoring is reduced.                 | <b>Single Model</b>         |
| 3. E. Fehr and L. Goette<br>"Do Workers Work More if Wages Are High? Evidence from a Randomized Field Experiment"<br><i>American Econ. Rev.</i> , 2007   | Random assignment of temporary increase in piece rate for bicycle messengers. Neoclassical model of intertemporal labor supply contrasted with reference-dependent preferences.         | <b>Competing Models</b>     |
| 4. P. Todd and K. Wolpin<br>"Assessing the Impact of a School Subsidy Program in Mexico: Using a Social Experiment to Validate a Dynamic Behavioral Model of Child Schooling and Fertility"<br><i>American Econ. Rev.</i> , 2006 | Random assignment of schooling subsidies. Village-level assignment. Dynamic structural model of fertility and schooling fit to control group and used to forecast experimental impacts. | <b>Parameter Estimation</b> |

Notes: Studies selected and summarized by authors. See text for description of relevant universe of studies, and classification system.

The fourth *Parameter Estimation* category includes studies that analyze field experiments using fully specified models. The estimation of the underlying parameters of the model allows for welfare and policy evaluations that are not possible otherwise. An interesting example is Todd and Wolpin (2006), who specify a dynamic choice model for schooling and fertility decisions of families in rural Mexico. They estimate the model parameters using data from the *control* group of the PROGRESA experiment, and then compare the predicted versus actual responses for the *treatment* group, who received financial incentives to participate in health, education, and nutrition programs.

## The Role of Theory in Experiments Since 1975

In this section, we turn to a quantitative analysis of the role that theory has played in field experiments published in five top journals over the past decades since 1975. To provide a useful contrast, we also classified all laboratory experiments published in five top journals, including laboratory-like experiments conducted in a field environment (labeled “artifactual field experiments” in Harrison and List, 2004).



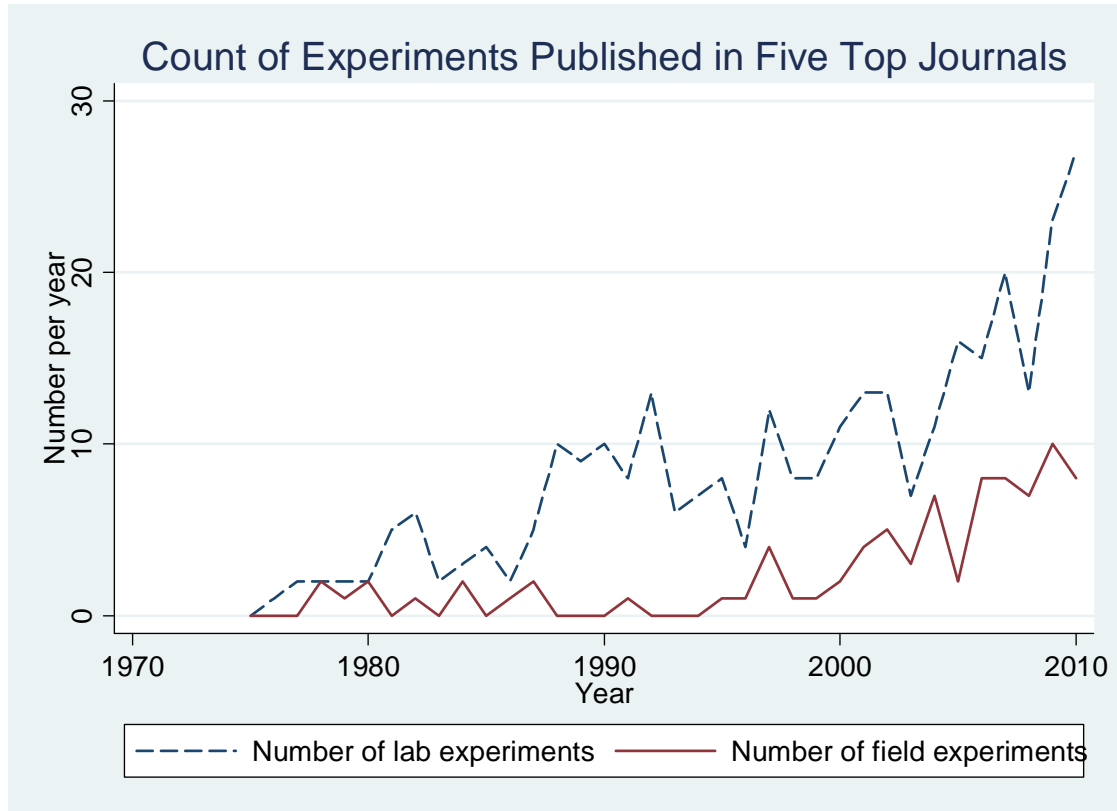
Figure 1 displays a count of all published field and laboratory experiments under these definitions. In addition, Table 2 lists all the field experiments classified, with the classification by the content of theory, as well as a rough categorization by field, and a measure of impact using a count of Google Scholar citations as of April 2011.

<Table 2 about here>

Until the mid-nineties the number of field experiments published in top journals was small. Between 1975 and 1984, eight field experiments were published in top-five journals, seven of which are analyses of the negative income tax experiments discussed later in this paper. Between 1985 and 1994, four more field experiments were published, including the Blank (1991) study of the impact of double anonymity in the refereeing process. Nearly all of these early field experiments are broadly in the area of labor economics, and they include several highly influential papers by the number of citations, including the LaLonde (1986) study of program evaluation methods (984 citations).

Since 1995, the number of field experiments has increased steadily, while the diversity of subject matter has also expanded to include such areas as behavioral economics (15 papers by our count), development economics (15 papers), public economics (13 papers, including the charity experiments), and industrial organization (8 papers, including the auction experiments). Since 1995, the authors with the most published field experiments by our categorization are John List (12 papers), Dean Karlan (5 papers), Esther Duflo (4 papers), and Joshua Angrist, Marianne Bertrand, Uri Gneezy, James Heckman, Lawrence Katz, Jeffrey Kling, Michael Kremer, Jeffery Liebman, and Sendhil Mullainathan (all with 3 papers each).

In the past six years the number of field experiments published has averaged 8-10 per year. Over our 36-year sample period, a total of 84 field experiments were published in the top-five journals.



**Figure 1. Number of laboratory and field experiments published in five top economics journals from 1975 to 2010**

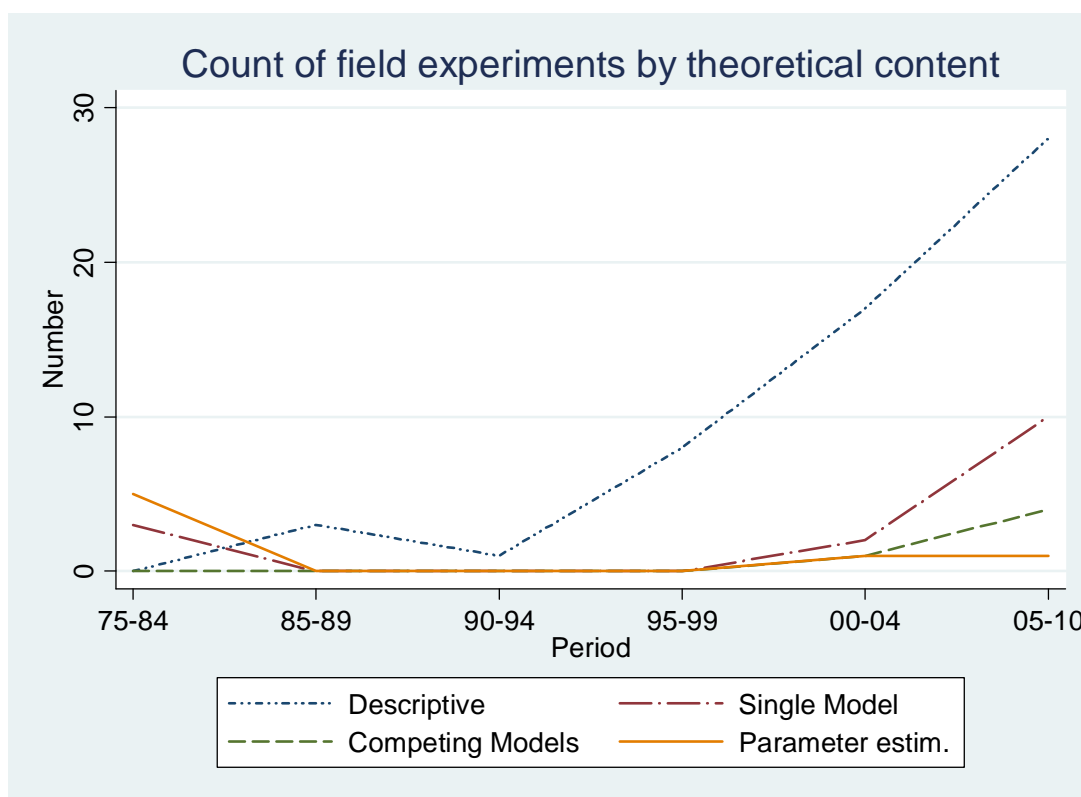
Compared to field experiments, laboratory experiments are far more common. In every year since 1981, more laboratory experiments than field experiments were published in the top-five journals. Between 1985 and 1995, the number per year ranged from five to ten, resulting in a total of 82 laboratory experiments, compared to only five field experiments in the same period. The flow of published laboratory experiments increased even further to 15-25 articles per year from 2005 to 2010. Indeed, in 2010, laboratory experiments account for 9.3 percent of all articles in these five top journals, compared to 2.5 percent for field experiments. The total number of laboratory experiments in our sample is 308, three and one-half times the number of field experiments.

The *American Economic Review* accounts for 54 percent of all laboratory studies in our data, followed by *Econometrica* (19 percent) and the *Quarterly Journal of Economics* (13 percent). Field experiments are more evenly distributed across journals, with the *American Economic Review* (35 percent) and the *Quarterly Journal of Economics* (27

percent) publishing the most, followed by *Econometrica* (19 percent). Within each of these journals, the trends over time are similar to the ones documented in Figure 1.

How many of these experiments fall into each of our four categories for the content of theory? Figure 2 shows the numbers in each category for field experiments for the initial decade of our sample period (1975-84), and for subsequent five-year periods (except for 2005-2010, which includes 6 years).

Interestingly, the field experiments published from 1975 to 1984 were all model-based: nearly all these papers used labor supply models to study the negative income tax experiments. The few field experiments published in 1985-89 and 1990-94 were all descriptive; so too were the eight field experiments published from 1995 to 1999. Among the 21 field experiments published in the 2000-04 period, 17 are descriptive while 4 have a higher theoretical content (as judged by our criteria): two with a *Single Model*, one with *Competing Models*, and one study with *Parameter Estimation*. The first field experiment with an explicit theoretical framework published in the post-1984 period is the Nagin et al. (2002) paper described above (in the *American Economic Review*). In the most recent 2005-2010 period theory has played a more important role in field experiments, with ten experiments with a *Single Model*, four with *Competing Models*, and one study with *Parameter Estimation*. Still, the dominant category remains *Descriptive*, with 28 articles.



**Figure 2. Field experiments by theoretical content**

Overall, 68 percent of the 86 field experiments published in top-five journals are *Descriptive*, 18 percent contain a *Single Model*, 6 percent contain *Competing Models*, and 8 percent of field experiments contain a model with *Parameter Estimation*.

The patterns are quite similar across journals, including *Econometrica* and the *Review of Economic Studies*. While empirical papers in these two journals are in general more likely to include models, in the case of field experiments the models are typically statistical, rather than economic models.

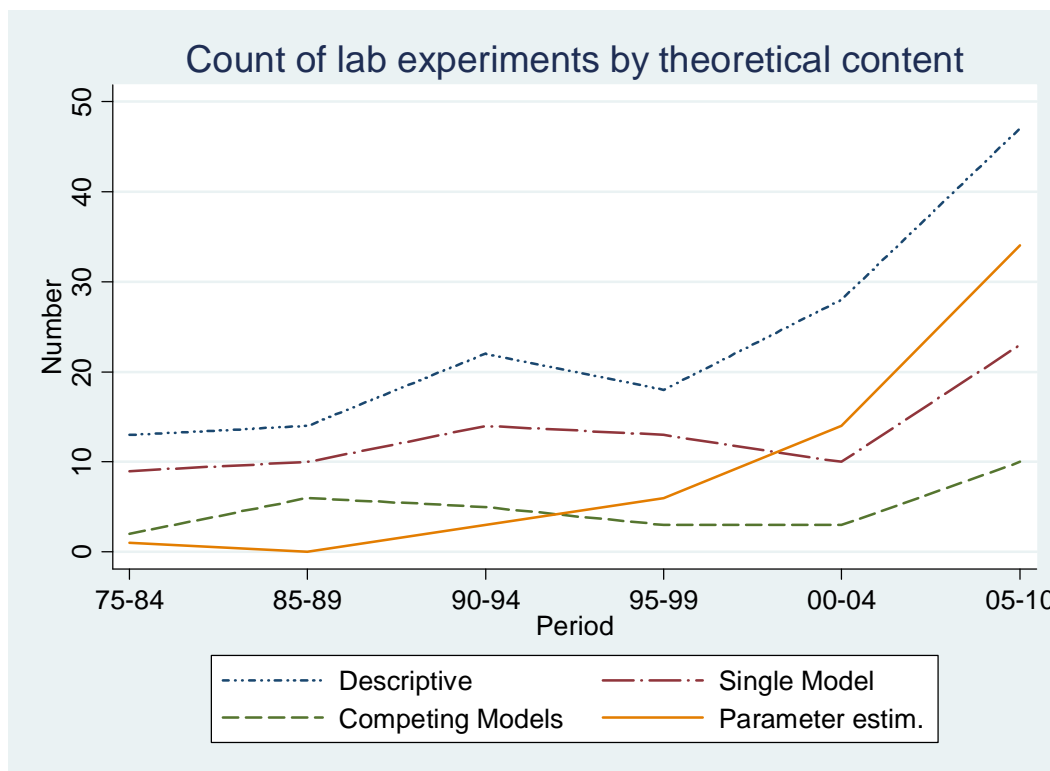
Using the citation counts, we evaluate the impact of the different categories of field experiment. In the period from 1995 to 2004, among papers with at least 200 citations, 16 out of 18 studies (89 percent) are *Descriptive*, which is in line with the share among all field experiments in those years (25 out of 29 studies). In the period from 2005 on, among papers with at least 100 citations, 8 of 13 (62 percent) are *Descriptive*, which is again in line with the overall share in this period (28 out of 43 studies). This evidence,

which is necessarily tentative given the small sample size, suggests that the citation-based measure of impact is similar across studies with different theoretical content.

Next, we consider the break-down by modeling content for laboratory experiments, as shown in Figure 3. The results are quite different. While the descriptive type of experiments has been, and remains, the most common type of laboratory experiment, model-based experiments (either with a single model or with competing models) have been relatively common since the 1970s. The main discernible trend in the last decade is an increase in the number of laboratory experiments with parameter estimation. This latter category includes, among others, the estimation of quantal response equilibria models, which provide a solution to game theory problems in a situation of bounded rationality; models of k-levels of thinking, in which the decisions of agents depend on how many levels of iteration they perform and they think other players will perform; and differing aspects of time and risk preferences.

Overall, it is clear that the role of explicit theoretical models is very different in laboratory than in field experiments: 26 percent of the laboratory experiments contain a *Single Model*, 9 percent contain *Competing Models*, and 19 percent of papers contain a model with *Parameter Estimation*, while only about one-half (46 percent) are *Descriptive* in nature.

These patterns differ by journal. In particular, *Econometrica* and the *Review of Economic Studies* have a higher incidence of model-based experiments than the other journals. In the last decade, the most common type of laboratory experiment in these two journals is one with *Parameter Estimation*.



**Figure 3. Laboratory experiments by theoretical content**

This brief historical review shows how different the role of theory is in laboratory and field experiments. Models have always played a key role in laboratory experiments, with an increasing trend. Field experiments have been largely descriptive, with only a recent increase in the role for models. In the two journals in our group of five typically most devoted to theory, *Econometrica* and the *Review of Economic Studies*, the most common laboratory experiment in the last decade is an experiment with a model and including *Parameter Estimation*, while the most common field experiment is descriptive.

The question then arises: What would be gained, and what would be lost, if field experiments were more like laboratory experiments, with respect to theory? We discuss this question using three exemplar types of field experiments: gift exchange experiments, charitable giving field experiments, and negative income tax studies.

### **Gift Exchange Field Experiments**

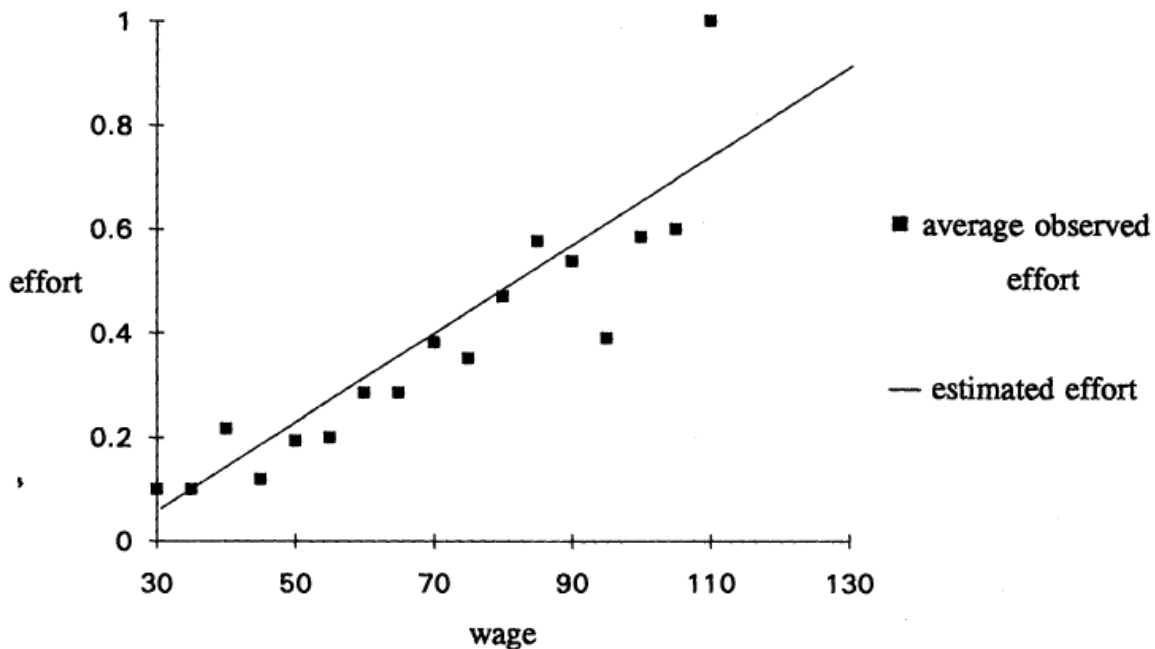
Akerlof (1982) argued that a gift exchange mechanism between employers and employees can play an important role in labor markets. If employees respond to a kind

wage offer by working harder, employers may find it optimal to offer wages above the reservation utility. Gift exchange, hence, is a possible rationale for efficiency wages.

This theory has proven hard to test empirically. For one thing, the repeated nature of employment contracts makes it difficult to separate genuine gift exchange from repeated game equilibria, in which the worker exerts extra effort in anticipation of future compensation and so on. In a genuine gift exchange, instead, the worker exerts extra effort because the “gift” by the employer induces pro-social behavior towards the employer.

In a highly-cited laboratory study, Fehr, Kirchsteiger and Riedl (1993) test for gift exchange. In the experiment, some subjects are assigned the role of firms, others the role of workers. Firms move first and make a wage offer  $w \in \{0, 5, 10, \dots\}$ . Workers then choose effort  $e \in [0.1, 1]$ . Workers and firms engage in one-shot interactions, so repeated-game effects are eliminated by design. Since effort is costly, the subgame perfect equilibrium strategy for self-interested workers is to exert the minimal effort  $e^* = 0.1$ , no matter what the wage offer. In anticipation of this, self-interested firms should offer workers their reservation utility, in this case  $w^* = 30$ .

Fehr, Kirchsteiger, and Riedl observe behavior that is starkly different from these predictions. Almost all subjects in the role of firms offer wages higher than 30, and subjects in the role of workers respond by exerting higher effort, as shown in Figure 4. This is precisely, in a laboratory setting, the gift exchange that Akerlof (1982) postulated. The reciprocal behavior of the workers makes it rational for firms to offer efficiency wages. A number of laboratory experiments have confirmed and extended the findings of this paper.



**Figure 4: Gift Exchange in a Laboratory Experiment**

**Source: Reproduced from Fehr, Kirchsteiger, and Riedl (1993).**

As interesting as this evidence is, one may argue that behavior in an actual employment contract differs from behavior in the laboratory. Yet, employment relationships with their repeated nature make testing of gift exchange behavior very difficult.

Gneezy and List (2006) designed a field experiment that resolves this difficulty. They hire workers for two tasks, coding library books and running a fund-raising drive. They make it clear that the jobs are one-time tasks, hence removing repeated-interaction incentives. Once subjects show up for their task, a sub-set is randomly assigned a surprise pay of \$20 per hour, while the control group is paid \$12 per hour, as promised. Gneezy and List then examine whether effort responds to the higher pay, as predicted by the gift exchange hypothesis. Notice that the higher pay is a fixed payment per hour, and as such does not alter the incentives to exert effort. The main finding in the paper is that work effort is substantially higher in the first three hours of the job in the gift treatment relative to the control treatment, but it is indistinguishable after that. This suggests that gift exchange is present, but short-lived. This innovative design spawned a whole literature of field experiments using similar short-term, but real, employment contracts.



What neither Gneezy and List (2006) nor most of the follow-up papers do is to provide a model for the observed behavior. As such, they are *Descriptive* field experiments. However, while gift exchange is indicative of non-standard preferences (else the worker would not reciprocate in a one-shot interaction), various models of social preferences can explain the evidence.

Two prominent classes of explanations are inequity aversion and reciprocity. Under inequity aversion, put forward by Fehr and Schmidt (1999) and Bolton and Ockenfels (2000), individuals dislike inequity: while individuals do want higher payoffs for themselves, they are willing to forgo some payoff to help another player who is behind them – though not someone who is ahead of them. This simple model of social preferences has been successful in accounting for behavior in a variety of contexts, including behavior in the dictator game, the ultimatum game, and gift exchange in the laboratory. In the Fehr et al. (1993) experiment, the “firm” falls behind by paying a (higher) wage. The worker can mitigate this inequity by exerting effort which benefits the firm with, at least initially, limited cost (since the cost function is convex). The model also predicts that the worker will not put this effort if the firm has not paid a generous wage. In this latter case, the firm is ahead in payoffs and putting in effort would increase, not decrease, inequality.

Under reciprocity models instead (such as the intention-based models in Rabin, 1993, and Dufwenberg and Kirchsteiger, 2004, or type-dependent preferences in Levine, 1998, or action-based models as in Cox, Friedman and Sadiraj, 2008), individuals have positive social preferences towards others who they think are nice or behave nicely, but not (as much) towards individuals who are not nice. Under these models, workers exert effort if the firm pays a higher wage in the laboratory gift exchange game because of the inference workers make about how nice the firm is. Conversely, they do not exert effort under a low wage because they do not care for firms that prove to be selfish.

Can gift exchange experiments in the field then help separate the two explanations? It is simple to show that they do, even though this point has not been made in the papers cited above. The inequity aversion model predicts gift exchange in the laboratory because the generous wage payment by the firm causes the firm to fall behind in payoffs relative to the worker, triggering the inequity-diminishing effort by the worker. But in the field experiment, it is highly implausible that a higher wage payment by the firm for a

six-hour task causes the firm to fall behind in payoffs relative to the workers. But if the “gift” payment does not alter the inequity between the worker and the firm, it will not induce gift exchange behavior. Hence, any observed gift exchange in firms cannot be due to inequity aversion but to other social preferences such as reciprocity. This point applies to other economic settings where gifts are given to influence behavior, such as gifts to doctors in the pharmaceutical industry (Malmendier and Schmidt, 2010), or vote-buying in the case of politicians (Finan and Schechter, 2010). These gift-exchange patterns cannot be explained by inequity aversion, but only by some of the existing reciprocity-based theories.

Adding a simple model of two (or more) competing social-preference models would thus add insights beyond the *Descriptive* contribution of the field experiments. Moreover, using a model of reciprocal preferences, one can ask how much reciprocity is implied by the observed gift exchange in the field. In Gneezy and List (2006), the increase in pay raises productivity in book coding (temporarily) by 30 percent. But did that gain require great effort, in which case it indicates substantial reciprocation, or only a minimal increase in effort, and thus not much reciprocation? Estimating the extent of reciprocity would require knowing the shape of the cost function of effort. This can be done by randomizing the piece rate. As such, additional experimental treatments can be designed to estimate the nuisance parameters (in this case the curvature of the cost of effort) and shed light on the parameters of interest (the extent of reciprocity).

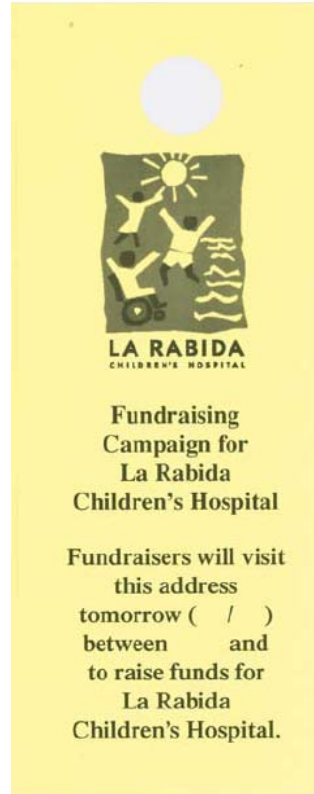
To summarize, the gift exchange experiments suggest that there is an important role played by both types of experimental evidence: The laboratory experiments in Fehr et al. (1993) were the first to suggest an experimental methodology to test for gift exchange, and found support for it in the laboratory. The Gneezy and List (2006) field experiment was a milestone in that it proposed a design for gift exchange in a real employment contract unconfounded by repeated game effects. While this field experiment falls in the *Descriptive* category, follow-up modeling can clarify its implications for the body of theory on social preferences. Furthermore, studies that structurally estimate these parameters could build on the design of Gneezy and List. Scientific progress is often achieved by a sequence of papers, each adding to the previous work.

## Charitable Giving Field Experiments

A series of field experiments have transformed the charitable giving field from an area mostly focused on modeling and stylized facts to one focused on experimental findings. A trail-blazing field experiment was List and Lucking-Reiley (2002). In a mailer requesting funds for a research center, the authors randomized both the seed money (the funding already available) and whether funds would be refunded in case the fund-raising targets were not met. This experiment was motivated by Andreoni's signalling model of charitable giving. However, since the List and Lucking-Reiley (2002) paper does not contain a model, we categorize it as *Descriptive*. Most recent field experiments in the area follow List and Lucking-Reiley: they are motivated by models on charitable giving, but they are ultimately *Descriptive* (for example, Falk, 2007).

In this Section, we discuss the role that theory played in a field experiment on charitable giving run by two of the authors of this paper, Stefano DellaVigna and Ulrike Malmendier, together with John List. The idea of the paper was to attempt to discriminate between two sets of reasons for giving to a charity when asked for a donation. One reason is that the act of giving is associated with a utility increase, whether due to altruism, warm glow, or prestige. Alternatively, individuals may actually *dislike* giving money to a charity but feel worse saying no to the solicitor. In this case charity giving is due to the social pressure that the individuals experience when being asked. These two motivations for giving have very different welfare implications for the giver: giving is welfare-increasing for the donor in the first case, but welfare-diminishing for the donor in the second case.

In the discussion of the experimental design, we settled on a door-to-door campaign where we would randomize the extent to which people are informed about the upcoming fund-raising campaign. In the treatment group, but not in the control group, we would post a flyer on the door-knob of the household, reproduced in Figure 5, informing them of the upcoming fund-raiser. Households could then vote with their feet—if giving is mostly due to altruism, households in the treatment group would sort into staying at home and give; if giving is mostly due to social pressure, they would sort out to avoid being asked.



**Figure 5. Example of the flyer in a charitable giving experiment**

**Source: Flyer used by DellaVigna, List, and Malmendier (2010).**

The initial plan for the field experiment was in the *Descriptive* line of previous work: we intended to test a hypothesis which was intuitively suggested by theory, but without actually making the underlying model explicit. After some discussion, though, we decided to write down a model to clarify what assumptions we were implicitly making. We assumed a cost function of shifting the probability of being at home (in response to the flyer), and we allowed for competing models to explain sorting and giving behavior: altruism on the one hand and a social pressure cost from turning down an in-person giving request on the other hand.

In our case, the dividends from writing the model were substantial. In addition to clarifying the assumptions needed (for example, that there is no social pressure cost from avoiding the solicitor by not answering the door), the model suggested novel predictions. One such prediction relates to the size of donations. In our model, social pressure drives small donations, but not larger ones. Hence, if social pressure is responsible for the ob-

served donations, the flyer treatment should lower small donations, but not larger ones. The model also suggested new treatments. In particular, we added an “Opt-Out” treatment in which the flyer includes a box that can be checked if the household does not “want to be disturbed.” This treatment makes sorting easier—that is, it lowers the cost of avoiding the solicitor relative to the regular flyer without opt-out box. Hence any (additional) decrease in giving allows us to identify social pressure more directly and to address confounding explanations such as information or self- and other-signaling models.

In summary, making the model explicit before running the experiment made for a tighter and more informative test of the initial hypothesis.

In addition, we realized that, were it not for one nuisance parameter, we would be able to estimate the key parameters of the model, including the social pressure cost of saying no to an in-person request, and the extent of altruism. The nuisance parameter is the elasticity of the cost of sorting in and out of the home, a key parameter to make inferences. Suppose for example that the flyer reduces the probability of home presence by 4 percentage points – should that be considered much or is little? Unfortunately, none of the experimental treatments allowed us to “monetize” the magnitude and estimate this elasticity parameter.

This led us to think of other ways to estimate this parameter. In the end, while still in the design stage, we decided to run a parallel field experiment specifically designed for the purpose. We posted flyers announcing that “Researchers will visit this address tomorrow ( / ) between ... and ... to conduct an X-minute survey. You will be paid \$Y for your participation.” Across treatments we varied the time duration X (5 or 10 minutes) and the payment Y (\$0, \$5, or \$10). The responsiveness in the presence at home with respect to the duration and the payment provided the identification to the elasticity parameters, hence allowing us to back out all other parameters. Indeed, in the end these survey treatments made up the bulk of our field experiment, even though their only purpose was to estimate a nuisance parameter.

The reduced-form results in DellaVigna, List, and Malmendier (2010) point to the importance of social pressure for solicited donations, with the most important piece of evidence being the fact that the flyer with opt-out option lowers donations significantly, and especially small donations. As discussed above, this is a key prediction of the social

pressure framework which we had not honed in until we wrote the model. As such, writing the model provided us with a tighter reduced-form test.

What do the survey treatments and the ensuing parameter estimation add to these results? We estimate the effect of a fund-raising campaign on the welfare of the households contacted. In a model with no social pressure, the welfare effect of a campaign can only be positive, since a donor can always costlessly say no. But in the presence of social pressure, this free-disposal condition does not hold: the benefits of a campaign for the willing donors have to be weighed against the cost non-donors pay for being asked and saying no, which we estimate to be about \$4 for a local charity. In addition to this cost for non-donors, we estimate that as many as 50 percent of the donors would have preferred not to be asked, because social pressure induces them to give when they would not have given otherwise, or give more than they otherwise would.

Taking into account these forces, our benchmark specification indicates that our door-to-door campaign induces a welfare loss of about \$1 on average per household contacted (including households that were not at home and hence did not suffer a welfare loss, and not counting the benefits associated with the public good provision). An interesting and counterintuitive result is that raising money for the local and well-liked favorite charity is associated with more negative welfare impacts than raising money for an out-of-state and lesser-known charity. More people are willing to donate to the local charity, but at the same time, the social pressure cost of saying “no” to the local charity is significantly higher, and the second force dominates. These latter findings, which of course require some parametric assumptions, complement the descriptive findings.

## **Negative Income Tax Experiments**

The two previous examples suggest that, in many experimental settings, much can be gained from a careful consideration of the predictions of economic models. But is it always advantageous to have a model with parameter estimation? In this Section we consider the case of the negative income tax experiments, one of the most famous large-scale social experiments conducted in the United States. Funded by the Office of Economic Opportunity from 1968 to 1972, this experiment was designed to test the effects of a neg-

ative income tax – a simplified two-parameter income support system proposed by Milton Friedman in the 1950s involving a guaranteed baseline amount of income (the first parameter), which is then phased out at a constant rate (the second parameter) as income is earned. Dozens of high-profile economists were involved in the design and analysis of this experiment.

The experimental design was closely tied to a specific parametric model: Rather than implement a simple “two-group” experimental design, the experiment included a total of eight different treatment arms, each with a specific value for the “guarantee level” (that is, the level of income support for a family with no earnings) and for the program tax rate. A complex optimal assignment model, developed by John Conlisk and Harold Watts, was designed to maximize the efficiency of the experiment, assuming a (parametric) model of the likely responses to the experimental incentives.

In principle, the design could have provided estimates of the incentive effects of various combinations of the guarantee level and tax rate. However, with the very small sample sizes (1,350 subjects, with 750 members of the treatment group, and 46-138 treatments per arm), even the pooled experimental impacts were quite imprecise. The only possible inferences that could be made from the data were under the assumption of the structural model.

Similarly complex designs were employed in the Rural Income Maintenance Experiment (operated in Iowa and North Carolina between 1969 and 1973), the Gary Income Maintenance Experiment (operated in Gary Indiana between 1971 and 1974), and the Seattle-Denver Income Maintenance Experiment (SIME-DIME), which ran between 1971 and 1982. As in the earlier negative income tax experiments, the SIME-DIME experiment was hampered by a small sample size – the SIME-DIME sample would have had to have been eight times larger to yield statistically significant treatment effect estimates for even the largest arm of the design.

From today’s perspective, the obvious comfort that analysts at the time had with a model-based assigned mechanism is surprising. Equally remarkable, perhaps, was the nearly universal adoption of model-based analysis methods for the negative income tax experiments (for example, see the analysis in Johnson and Pencavel, 1982). As pointed

out by Ashenfelter and Plant (1990), the final report of the SIME-DIME experiment did not include any “non-parametric” estimates of the impact of treatment.

As a result of the frustrations in dealing with the complex designs of the negative income tax experiments (and with the confusing message that emerged from such designs) many respected analysts adopted the view that social experiments should be designed as simply as possible. For example, Hausman and Wise (1985, p. 188) argued: “[W]e propose as a guiding principle the experiments should have as a first priority the precise estimation of a single or a small number of treatment effects.” Subsequent social experiments – particularly those that focus on new programs –have tended to follow this advice. As noted by Greenberg, Shroder and Onstott (1999) in this journal, 80 percent of the social experiments initiated after 1983 had only a single treatment-control contrast. This shift away from designs that explicitly attempt to model response variation to multiple treatments and toward a single manipulation has led to a new round of criticism that the social experiments are often “black boxes” that “... contribute next to nothing to the cumulative body of social science knowledge...” (Heckman and Smith, 1995, p. 108).

## **Conclusions**

Over the last two decades, economics has witnessed a dramatic expansion of experimental research. Both laboratory and field experiments share the common advantage of studying a controlled setting in order to evaluate treatment effects. There is, however, as we documented, a noticeable difference in the evolution of these two types of experimental research: Laboratory experiments feature a much closer link to theory than field experiments.

Examples from studies of gift exchange and charitable giving illustrate that, while we can certainly learn from descriptive studies, developing a fully specified behavioral model and obtaining estimates of the key parameters from that model can provide additional insights. This process can follow from models and estimates that are obtained in follow-up papers, as may happen for the gift exchange experiments, or could be part of the design of the initial field experiment, as in the charity experiment described above. In



this way, field experiments need not differ from laboratory experiments with respect to the guiding role that theory can play in testing hypothesis on behavior.

The negative income tax experiment, on the other hand, makes it clear that there is no simple answer as to the optimal role of modeling in field experiment. Reliance on a model is not always a plus, particularly in the evaluation of complex social programs that may affect a range of behaviors through multiple channels.

## Acknowledgements

We thank for helpful comments the editors, as well as Oriana Bandiera, Iwan Barankay, Glenn Harrison, Matthew Rabin, David Reiley, and participants in Berkeley, at the Wharton Conference on Field Experiment, and at the 2011 ASSA conference in Denver. We thank Ivan Balbuzanov, Xiaoyu Xia, and a very dedicated group of undergraduate students for excellent research assistance.

## References

Akerlof, George A. "Labor Contracts as Partial Gift Exchange." *The Quarterly Journal of Economics*. Vol. 97, No. 4: 543-569

Joshua D. Angrist. "Lifetime Earnings and the Vietnam Era Draft Lottery: Evidence from Social Security Administrative Records" *American Economic Review* Vol. 80, No. 3 (Jun., 1990), pp. 313-336.

Ashenfelter, Orley and Mark W. Plant. "Nonparametric Estimates of the Labor Supply Effects of Negative Income Tax Programs. *Journal of Labor Economics* 8 (January 1990, part 2), pp. S396-S415.

Bandiera, Oriana , Iwan Barankay,. and Imran Rasul. 2009. "Social Connections and Incentives in the Workplace: Evidence From Personnel Data." *Econometrica*, 77: 1047–1094.

Bandiera, Oriana , Iwan Barankay,. and Imran Rasul. 2007, "Incentives for Managers and Inequality Among Workers: Evidence from a Firm Level Experiment", *Quarterly Journal of Economics* 122: 729-74.

Barberis, Nicholas, and Wei Xiong. 2009. "What Drives the Disposition Effect? An Analysis of a Long-Standing Preference-Based Explanation." *Journal of Finance*, Vol. 64, Issue 2, pages 751–784.

Bolton, Gary E and Axel Ockenfels. Mar., 2000. "ERC: A Theory of Equity, Reciprocity, and Competition." *The American Economic Review*, Vol. 90, 166–193.

Gary S. Becker. "Irrational Behavior and Economic Theory" *Journal of Political Economy*, Vol. 70, No. 1 (Feb., 1962), pp. 1-13

Rebecca M. Blank. "The Effects of Double-Blind versus Single-Blind Reviewing: Experimental Evidence from The American Economic Review" *American Economic Review*, Vol. 81, No. 5 (Dec., 1991), pp. 1041-1067

Card, David. and Dean. R. Hyslop. 2005. "Estimating the Effects of a Time-Limited Earnings Subsidy for Welfare-Leavers." *Econometrica*, 73: 1723–1770

Chetty, Raj, Adam Looney and Kory Kroft. 2009. "Salience and Taxation: Theory and Evidence." *The American Economic Review*. Vol.99(4) : 1145-1177

Cox, James C; Friedman, Daniel and Vjollca Sadiraj. 2008. "Revealed Altruism." *Econometrica*, Vol. 76, 31–69.

Angus Deaton. "Understanding the Mechanisms of Economic Development" *Journal of Economic Perspectives*, vol. 24, iss. 3, pp. 3-16, August 2010.

DellaVigna, Stefano, Ulrike Malmendier, and John. A. List. 2010. "Testing for Altruism and Social Pressure in Charitable Giving" Working paper.

Falk, Armin. 2007. "Gift Exchange in the Field". *Econometrica*, 75: 1501–1511.

Fehr, Ernst and Lorenz Goette. 2007. "Do Workers Work More if Wages Are High? Evidence from a Randomized Field Experiment." *The American Economic Review*. Vol. 97, No. 1: 298-317

Fehr, Ernst , Georg Kirchsteiger and Arno Riedl. 1993. "Does Fairness Prevent Market Clearing? An Experimental Investigation." *The Quarterly Journal of Economics*. Vol. 108, No. 2 : 437-459

Fehr, Ernst, and Klaus M. Schmidt. 1999. "A Theory of Fairness, Competition, and Cooperation." *Quarterly Journal of Economics*, 114(3): 817–68.

Finan, Fred and Laura Schechter. 2010. "Vote-buying and Reciprocity," Working paper.

Forsetlund Louise, Iain Chalmers and Arild Bjørndal. 2007. "When Was Random Allocation First Used To Generate Comparison Groups In Experiments To Assess The Effects Of Social Interventions?" *Economics of Innovation and New Technology*. Vol. 16(5) : 371-384

Gneezy, Uri. and List, John. A. 2006. "Putting Behavioral Economics to Work: Testing for Gift Exchange in Labor Markets Using Field Experiments". *Econometrica*, 74: 1365–1384.

Harrison, Glenn W., and John A. List. 2004. "Field Experiments." *Journal of Economic Literature*. 42 : 1013–1059.

Greenberg, David, Mark Shroder and Matthew Onstott. "The Social Experiment Market." *Journal of Economic Perspectives* 13 (Summer 1999), pp. 157-172.

Hausman, Jerry A. and David A. Wise. "Technical Problems in Social Experimentation: Cost versus Ease of Analysis." In Jerry A Hausman and David A. Wise, editors. *Social Experimentation*. Chicago: University of Chicago Press, 1985.

Heckman, James J. and Jeffrey A. Smith. "Assessing the Case for Social Experiments." *Journal of Economic Perspectives* 9 (Spring 1995), pp. 85-110.

Johnson, Terry R. and Pencavel, John H. 1982. "Forecasting the Effects of a Negative Income Tax Program". *Industrial and Labor Relations Review*. Vol. 35, No. 2 : 221-234

Kling, Jeffrey R., Jeffrey B Liebman and Lawrence F Katz. 2007. "Experimental Analysis of Neighborhood Effects." *Econometrica*, 75: 83–119.

LaLonde, Robert J. 1986. "Evaluating the Econometric Evaluations of Training Programs with Experimental Data." *The American Economic Review*. Vol. 76, No. 4 : 604-620

Levine, David K. 1998. "Modeling Altruism and Spitefulness in Experiments." *Review of Economic Dynamics*, Vol. 1, 593–622.

List, John A. "Does Market Experience Eliminate Market Anomalies?" *Quarterly Journal of Economics*, 2003, Vol. 118, No. 1 : 41-71

List, John A. and David Lucking-Reiley. 2002. "The Effects of Seed Money and Refunds on Charitable Giving: Experimental Evidence from a University Capital Campaign". *The Journal of Political Economy* Vol. 110, No. 1 : 215-233

Malmendier, Ulrike and Klaus Schmidt. 2010. "You Owe Me." Unpublished Manuscript.

Miguel, Edward and Michael Kremer. 2004. "Worms: Identifying Impacts on Education and Health in the Presence of Treatment Externalities." *Econometrica*, 72: 159–217.

Nagin, Daniel S., James B. Rebitzer, Seth Sanders and Lowell J. Taylor. 2002. "Monitoring, Motivation, and Management: The Determinants of Opportunistic Behavior in a Field Experiment." *The American Economic Review*. Vol. 92, No. 4 : 850-873

Odean, Terrance. 1998. "Are Investors Reluctant to Realize Their Losses?" *Journal of Finance*, 53(5): 1775–98.

Rabin, 2010. "Improving Theory with Experiments, Improving Experiments with Theory, all so as to Improve Traditional Economics." Unpublished Manuscript.

Sacerdote, Bruce. 2001. "Peer Effects with Random Assignment: Results for Dartmouth Roommates." *Quarterly Journal of Economics*. Vol.116:2, 681-704

Todd, Petra E. and ; Kenneth I. Wolpin. 2006. "Assessing the Impact of a School Subsidy Program in Mexico: Using a Social Experiment to Validate a Dynamic Behavioral Model of Child Schooling and Fertility." *The American Economic Review*. Vol.96 (5): 1384-1417

**Table 2, Panel A. List of all Field Experiments Published in Top-5 Journals from 1975 to 2005**

| Year | Month | Journal | Pages     | Authors   | Title   | Classification   | Google     |             |
|------|-------|---------|-----------|---|---|------------------|------------|-------------|
|      |       |         |           |   |   |                  | Sch. Cites | Field       |
| 1978 | 12    | JPE     | 1103-1130 | Burtless, Gary and Jerry A. Hausman                   | The Effect of Taxation on Labor Supply: Evaluating the Gary Negative Income Tax Experiment        | Single Model     | 310        | Labor       |
| 1978 | 12    | AER     | 873-887   | Keeley, Michael C., Robins, Philip K., Spiegelman, F  | The Estimation of Labor Supply Models Using Experimental Data                                     | Parameter Estim. | 65         | Labor       |
| 1979 | 3     | EMA     | 455-473   | Hausman, Jerry A. and David A. Wise                   | Attrition Bias in Experimental and Panel Data: The Gary Income Maintenance Experiment             | Single Model     | 285        | Labor       |
| 1980 | 1     | RES     | 75-96     | Hausman, Jerry A. and David A. Wise                   | Discontinuous Budget Constraints and Estimation: The Demand for Housing                           | Parameter Estim. | 39         | Labor       |
| 1980 | 5     | EMA     | 1031-1052 | Tuma, Nancy B. and Philip K. Robins                   | A Dynamic Model of Employment Behavior: An Application to the Seattle and Denver Income Main      | Single Model     | 33         | Labor       |
| 1982 | 6     | AER     | 488-497   | Burtless, Gary and David Greenberg                    | Inferences Concerning Labor Supply Behavior Based on Limited-Duration Experiments                 | Parameter Estim. | 25         | Labor       |
| 1984 | 3     | EMA     | 363-390   | T. R. Johnson, J. H. Pencavel                         | Dynamic Hours of Work Functions for Husbands, Wives, and Single Females                           | Parameter Estim. | 62         | Labor       |
| 1984 | 9     | AER     | 673-684   | Plant, Mark W.  | An Empirical Analysis of Welfare Dependence   | Parameter Estim. | 43         | Labor       |
| 1986 | 9     | AER     | 604-620   | LaLonde, Robert J.                                    | Evaluating the Econometric Evaluations of Training Programs with Experimental Data                | Descriptive      | 984        | Labor       |
| 1987 | 6     | AER     | 251-277   | Manning, Willard G., Newhouse, Joseph P., Duan, J     | Health Insurance and the Demand for Medical Care: Evidence from a Randomized Experiment           | Descriptive      | 1165       | Health      |
| 1987 | 9     | AER     | 513-530   | Woodbury, Stephen A. and Robert G. Spiegelman         | Bonuses to Workers and Employers to Reduce Unemployment: Randomized Trials in Illinois            | Descriptive      | 154        | Labor       |
| 1991 | 12    | AER     | 1041-1067 | Blank, Rebecca M.                                     | The Effects of Double-Blind versus Single-Blind Reviewing: Experimental Evidence from The Ame     | Descriptive      | 222        | Labor       |
| 1995 | 6     | AER     | 304-321   | Ayres, Ian and Peter Siegelman                        | Race and Gender Discrimination in Bargaining for a New Car  | Descriptive      | 255        | IO          |
| 1996 | 1     | EMA     | 175-205   | Ham, John C. and Robert J. Lalonde                    | The Effect of Sample Selection and Initial Conditions in Duration Models: Evidence from Experime  | Descriptive      | 271        | Labor       |
| 1997 | 10    | RES     | 487-535   | Heckman, James J., Smith, Jeffrey and Nancy Cler      | Making the Most Out of Programme Evaluations and Social Experiments: Accounting for Heteroge      | Descriptive      | 362        | Labor       |
| 1997 | 10    | RES     | 537-553   | Manski, Charles F.                                    | The Mixing Problem in Programme Evaluation  | Descriptive      | 50         | Labor       |
| 1997 | 10    | RES     | 605-654   | Heckman, James J., Ichimura, Hidehiko and Petra       | Matching as an Econometric Evaluation Estimator: Evidence from Evaluating a Job Training Prog     | Descriptive      | 1740       | Labor       |
| 1997 | 10    | RES     | 655-682   | Eberwein, Curtis, Ham, John C. and Robert J. Lalor    | The Impact of Being Offered and Receiving Classroom Training on the Employment Histories of D     | Descriptive      | 98         | Labor       |
| 1998 | 6     | JPE     | 457-482   | Camerer, Colin F.                                     | Can Asset Markets Be Manipulated? A Field Experiment with Racetrack Betting                       | Descriptive      | 100        | Asset Pr.   |
| 1999 | 5     | QJE     | 497-532   | Krueger, Alan B.                                      | Experimental Estimates of Education Production Functions  | Descriptive      | 871        | Labor       |
| 2000 | 5     | QJE     | 651-694   | Heckman, James, Hohmann, Neil, Smith, Jeffrey ar      | Substitution and Dropout Bias in Social Experiments: A Study of an Influential Social Experiment  | Descriptive      | 89         | Labor       |
| 2000 | 9     | AER     | 961-972   | List, John A. and David Lucking-Reiley                | Demand Reduction in Multiunit Auctions: Evidence from a SportsCard Field Experiment               | Descriptive      | 177        | IO          |
| 2001 | 5     | QJE     | 607-654   | Katz, Lawrence F., Kling, Jeffrey R. and Jeffrey B. I | Moving to Opportunity in Boston: Early Results of a Randomized Mobility Experiment                | Descriptive      | 572        | Public      |
| 2001 | 5     | QJE     | 655-679   | Ludwig, Jens, Duncan, Greg J. and Paul Hirschfield    | Urban Poverty and Juvenile Crime: Evidence from a Randomized Housing-Mobility Experiment          | Descriptive      | 355        | Public      |
| 2001 | 7     | EMA     | 1099-1111 | Philipson, Tomas                                      | Data Markets, Missing Data, and Incentive Pay   | Descriptive      | 17         | Sv. Methods |
| 2001 | 12    | AER     | 1498-1507 | List, John A.   | Do Explicit Warnings Eliminate the Hypothetical Bias in Elicitation Procedures? Evidence from Fie | Descriptive      | 226        | Behavioral  |
| 2002 | 1     | EMA     | 91-117    | Abadie, Alberto, Angrist, Joshua and Guido Imbens     | Instrumental Variables Estimates of the Effect of Subsidized Training on the Quantiles of Trainee | Descriptive      | 217        | Labor       |
| 2002 | 2     | JPE     | 215-233   | List, John A. and David Lucking-Reiley                | The Effects of Seed Money and Refunds on Charitable Giving: Experimental Evidence from a Uni      | Descriptive      | 202        | Public      |
| 2002 | 9     | AER     | 850-873   | Nagin, Daniel S., Rebitzer, James B., Sanders, Set    | Monitoring, Motivation, and Management: The Determinants of Opportunistic Behavior in a Field E   | Single Model     | 159        | Labor       |
| 2002 | 12    | AER     | 1535-1558 | Angrist, Joshua, Bettinger, Eric, Bloom, Erik, King,  | Vouchers for Private Schooling in Colombia: Evidence from a Randomized Natural Experiment         | Descriptive      | 341        | Labor       |
| 2002 | 12    | AER     | 1636-1643 | List, John A.   | Preference Reversals of a Different Kind: The 'More Is Less' Phenomenon                           | Descriptive      | 71         | Behavioral  |
| 2003 | 2     | QJE     | 41-71     | List, John A.   | Does Market Experience Eliminate Market Anomalies?  | Descriptive      | 381        | Behavioral  |
| 2003 | 6     | JPE     | 530-554   | Grogger, Jeffrey and Charles Michalopoulos            | Welfare Dynamics under Time Limits  | Descriptive      | 105        | Public      |
| 2003 | 8     | QJE     | 815-842   | Duflo, Esther and Emmanuel Saez                       | The Role of Information and Social Interactions in Retirement Plan Decisions: Evidence from a Ra  | Descriptive      | 366        | Public      |
| 2004 | 1     | EMA     | 159-217   | Miguel, Edward and Michael Kremer                     | Worms: Identifying Impacts on Education and Health in the Presence of Treatment Externalities     | Descriptive      | 537        | Development |
| 2004 | 2     | QJE     | 49-89     | List, John A.   | The Nature and Extent of Discrimination in the Marketplace: Evidence from the Field               | Descriptive      | 86         | IO          |
| 2004 | 3     | EMA     | 615-625   | List, John A.   | Neoclassical Theory versus Prospect Theory: Evidence from the Marketplace                         | Competing Model  | 226        | Behavioral  |
| 2004 | 4     | RES     | 513-534   | Shearer, Bruce  | Piece Rates, Fixed Wages and Incentives: Evidence from a Field Experiment                         | Parameter Estim. | 123        | Labor       |
| 2004 | 9     | EMA     | 1409-1443 | Chattopadhyay, Raghabendra and Esther Duflo           | Women as Policy Makers: Evidence from a Randomized Policy Experiment in India                     | Single Model     | 286        | Development |
| 2004 | 9     | AER     | 991-1013  | Bertrand, Marianne and Sendhil Mullainathan           | Are Emily and Greg More Employable than Lakisha and Jamal? A Field Experiment on Labor Mar        | Descriptive      | 806        | Labor       |
| 2004 | 12    | AER     | 1717-1722 | Frey, Bruno S. and Stephan Meier                      | Social Comparisons and Pro-Social Behavior: Testing 'Conditional Cooperation' in a Field Experim  | Descriptive      | 226        | Behavioral  |
| 2005 | 2     | QJE     | 87-130    | Kling, Jeffrey R., Ludwig, Jens and Lawrence F. Ka    | Neighborhood Effects on Crime for Female and Male Youth: Evidence From a Randomized Housi         | Descriptive      | 222        | Public      |
| 2005 | 11    | EMA     | 1723-1770 | Card, David and R. Hyslop                             | Estimating the Effects of a Time-Limited Earnings Subsidy for Welfare-Leavers                     | Descriptive      | 77         | Labor       |



**Table 2, Panel B. List of all Field Experiments Published in Top-5 Journals from 2006 to 2010**

| Year | Month | Journal | Pages     | Authors  | Title   | Classification   | Google     |             |
|------|-------|---------|-----------|--|---|------------------|------------|-------------|
|      |       |         |           |  |   |                  | Sch. Cites | Field       |
| 2006 | 2     | JPE     | 1-37      | List, John A.  | The Behavioralist Meets the Market: Measuring Social Preferences and Reputation Effects in Act  | Descriptive      | 139        | Behavioral  |
| 2006 | 5     | QJE     | 635-672   | Ashraf, Nava, Karlan, Dean S. and Wesley Yin         | Tying Odysseus to the Mast: Evidence from a Commitment Savings Product in the Philippines       | Descriptive      | 289        | Development |
| 2006 | 5     | QJE     | 673-697   | Raymond Fisman, Sheena S. Iyengar, Emir Kamen        | Gender Differences in Mate Selection: Evidence from a Speed Dating Experiment                   | Single Model     | 92         | Labor       |
| 2006 | 5     | QJE     | 747-782   | Landry, Craig E., Lange, Andreas, List, John A., Pri | Toward an Understanding of the Economics of Charity: Evidence from a Field Experiment           | Single Model     | 115        | Public      |
| 2006 | 9     | EMA     | 1365-1384 | Gneezy, Uri and John A. List                         | Putting Behavioral Economics to Work: Testing for Gift Exchange in Labor Markets Using Field Ex | Descriptive      | 91         | Behavioral  |
| 2006 | 9     | AER     | 988-1012  | Bitler, Marianne P., Gelbach, Jonah B. and Hilary W  | What Mean Impacts Miss: Distributional Effects of Welfare Reform Experiments                    | Descriptive      | 102        | Public      |
| 2006 | 11    | QJE     | 1311-1346 | Duflo, Esther, Gale, William, Liebman, Jeffrey, Ors  | Saving Incentives for Low- and Middle-Income Families: Evidence from a Field Experiment with H  | Descriptive      | 107        | Public      |
| 2006 | 12    | AER     | 1384-1417 | Todd, Petra E. and Kenneth I. Wolpin                 | Assessing the Impact of a School Subsidy Program in Mexico: Using a Social Experiment to Valid  | Parameter Estim. | 110        | Development |
| 2007 | 1     | EMA     | 83-119    | Kling, Jeffrey R., Liebman, Jeffrey B. and Lawrence  | Experimental Analysis of Neighborhood Effects   | Descriptive      | 424        | Labor       |
| 2007 | 3     | AER     | 298-317   | Fehr, Ernst and Lorenz Goette                        | Do Workers Work More if Wages Are High? Evidence from a Randomized Field Experiment             | Competing Models | 165        | Behavioral  |
| 2007 | 4     | JPE     | 200-249   | Olken, Benjamin A.                                   | Monitoring Corruption: Evidence from a Field Experiment in Indonesia                            | Descriptive      | 321        | Development |
| 2007 | 8     | QJE     | 1007-1065 | Kremer, Michael and Edward Miguel                    | The Illusion of Sustainability  | Single Model     | 129        | Development |
| 2007 | 8     | QJE     | 1235-1264 | Banerjee, Abhijit V., Cole, Shawn, Duflo, Esther and | Remedying Education: Evidence from Two Randomized Experiments in India                          | Descriptive      | 211        | Development |
| 2007 | 9     | EMA     | 1501-1511 | Falk, Armin  | Gift Exchange in the Field  | Descriptive      | 80         | Behavioral  |
| 2007 | 11    | QJE     | 1639-1676 | Bertrand, Marianne, Djankov, Simeon, Hanna, Rem      | Obtaining a Driver's License in India: An Experimental Approach to Studying Corruption          | Descriptive      | 53         | Development |
| 2007 | 12    | AER     | 1774-1793 | Karlan, Dean S. and John A. List                     | Does Price Matter in Charitable Giving? Evidence from a Large-Scale Natural Field Experiment    | Descriptive      | 88         | Public      |
| 2008 | 5     | EMA     | 643-660   | Graham, Bryan S.                                     | Identifying Social Interactions through Conditional Variance Restrictions                       | Descriptive      | 61         | Labor       |
| 2008 | 6     | AER     | 1040-1068 | Karlan, Dean S. and Jonathan Zinman                  | Credit Elasticities in Less-Developed Economies: Implications for Microfinance                  | Descriptive      | 73         | Development |
| 2008 | 9     | AER     | 1553-1577 | Jensen, Robert T. and Nolan H. Miller                | Giffen Behavior and Subsistence Consumption   | Single Model     | 25         | IO          |
| 2008 | 11    | QJE     | 1329-1372 | de Mel, Suresh, McKenzie, David and Christopher      | Returns to Capital in Microenterprises: Evidence from a Field Experiment                        | Single Model     | 85         | IO          |
| 2008 | 11    | QJE     | 1373-1414 | Hastings, Justine S. and Jeffrey M. Weinstein        | Information, School Choice, and Academic Achievement: Evidence from Two Experiments             | Descriptive      | 51         | Public      |
| 2008 | 12    | AER     | 1829-1863 | Thornton, Rebecca L.                                 | The Demand for, and Impact of, Learning HIV Status  | Descriptive      | 52         | Development |
| 2008 | 12    | AER     | 1887-1921 | Schochet, Peter Z, Burghardt, John and Sheena M      | Does Job Corps Work? Impact Findings from the National Job Corps Study                          | Descriptive      | 37         | Public      |
| 2009 | 3     | AER     | 486-508   | Angelucci, Manuela and Giacomo De Giorgi             | Indirect Effects of an Aid Program: How Do Cash Transfers Affect Ineligibles' Consumption?      | Single Model     | 44         | Development |
| 2009 | 4     | RES     | 451-469   | Ariely, Dan, Gneezy, Uri, Loewenstein, George and    | Large Stakes and Big Mistakes   | Descriptive      | 70         | Behavioral  |
| 2009 | 5     | QJE     | 735-769   | Björkman, Martina and Jakob Svensson                 | Power to the People: Evidence from a Randomized Field Experiment on Community-Based Monit       | Descriptive      | 39         | Development |
| 2009 | 5     | EMA     | 909-931   | Charness, Gary and Uri Gneezy                        | Incentives to Exercise  | Descriptive      | 27         | Behavioral  |
| 2009 | 6     | JPE     | 453-503   | Bobonis, Gustavo J.                                  | Is the Allocation of Resources within the Household Efficient? New Evidence from a Randomized   | Competing Models | 18         | Labor       |
| 2009 | 6     | AER     | 864-882   | Cai, Hongbin, Chen, Yuyu and Hanming Fang            | Observational Learning: Evidence from a Randomized Natural Field Experiment                     | Descriptive      | 29         | Behavioral  |
| 2009 | 7     | RES     | 1071-1102 | Lee, David S.  | Training, Wages, and Sample Selection: Estimating Sharp Bounds on Treatment Effects             | Descriptive      | 79         | Labor       |
| 2009 | 9     | AER     | 1384-1414 | Angrist, Joshua D. and Victor Lavy                   | The Effect of High Stakes High School Achievement Awards: Evidence from a School-Centered R     | Descriptive      | 67         | Labor       |
| 2009 | 11    | QJE     | 1815-1851 | Leider, Stephen, Mobius, Markus M., Rosenblat, T     | Directed Altruism and Enforced Reciprocity in Social Networks: How Much is a Friend Worth?      | Descriptive      | 10         | Behavioral  |
| 2009 | 11    | EMA     | 1993-2008 | Karlan, Dean S. and Jonathon Zinman                  | Observing Unobservables: Identifying Information Asymmetries With a Consumer Credit Field Exp   | Competing Models | 129        | IO          |
| 2010 | 2     | QJE     | 1-45      | Cohen, Jessica and Pascaline Dupas                   | Free Distribution or Cost-Sharing? Evidence from a Randomized Malaria Prevention Experiment     | Single Model     | 41         | Development |
| 2010 | 2     | QJE     | 263-305   | Bertrand, Marianne, Karlan, Dean, Mullainathan, S    | What's Advertising Content Worth? Evidence from a Consumer Credit Marketing Field Experiment    | Single Model     | 51         | Behavioral  |
| 2010 | 4     | JPE     | 274-299   | Levav, Jonathan, Heitmann, Mark, Herrmann, Andr      | Order in Product Customization Decisions: Evidence from Field Experiments                       | Descriptive      | 10         | IO          |
| 2010 | 5     | QJE     | 515-548   | Jensen, Robert                                       | The (Perceived) Returns to Education and the Demand for Schooling                               | Descriptive      | 37         | Development |
| 2010 | 5     | QJE     | 729-765   | Anderson, Eric T. and Duncan I. Simester             | Price Stickiness and Customer Antagonism  | Descriptive      | 15         | IO          |
| 2010 | 6     | AER     | 958-83    | Landry, Craig E., Lange, Andreas, List, John A., Pri | Is a Donor in Hand Better Than Two in the Bush? Evidence from a Natural Field Experiment        | Single Model     | 5          | Public      |
| 2010 | 9     | AER     | 1358-98   | Chen, Yan, Harper, F. Maxwell, Konstan, Joseph a     | Social Comparisons and Contributions to Online Communities: A Field Experiment on MovieLens     | Single Model     | 22         | Behavioral  |
| 2010 | 12    | AER     | 2383-2413 | Ashraf, Nava, Berry, James and Jesse M. Shapiro      | Can Higher Prices Stimulate Product Use? Evidence from a Field Experiment in Zambia             | Competing Models | 51         | Development |

**Notes:** List of all papers published in top-5 journals from 1975 to 2010 which we classify as field experiments. For the categorization into 4 types by the role of theory, see text. The Google Scholar cite count is as of April 2011.